

## Aberystwyth University

### *DGIG-Net*

Liu, Xiyao; Ji, Zhong; Pang, Yanwei; Han, Jungong; Li, Xuelong

*Published in:*  
IEEE Transactions on Cybernetics

*DOI:*  
[10.1109/TCYB.2021.3049537](https://doi.org/10.1109/TCYB.2021.3049537)

*Publication date:*  
2022

*Citation for published version (APA):*

Liu, X., Ji, Z., Pang, Y., Han, J., & Li, X. (2022). DGIG-Net: Dynamic Graph-in-Graph Networks for Few-Shot Human-Object Interaction. *IEEE Transactions on Cybernetics*, 52(8), 7852-7864.  
<https://doi.org/10.1109/TCYB.2021.3049537>

#### **General rights**

Copyright and moral rights for the publications made accessible in the Aberystwyth Research Portal (the Institutional Repository) are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Aberystwyth Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Aberystwyth Research Portal

#### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

tel: +44 1970 62 2400  
email: [is@aber.ac.uk](mailto:is@aber.ac.uk)



## DGIG-Net: Dynamic Graph-In-Graph Networks for Few-Shot Human-Object Interaction

Journal:	<i>IEEE Transactions on Cybernetics</i>
Manuscript ID	CYB-E-2020-06-1555
Manuscript Type:	Regular Paper
Date Submitted by the Author:	29-Jun-2020
Complete List of Authors:	Liu, Xiyao; Tianjin University, School of Electrical and Information Engineering Ji, Zhong; Tianjin University, School of Electrical and Information Engineering Pang, Yanwei; Tianjin University, School of Electronic Information Engineering; Han, Jungong; University of Warwick, Data Science group Li, Xuelong; Northwestern Polytechnical University, School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL)
Keywords:	Few-Shot Learning, Human-Object Interaction, Meta-Learning, Dynamic Graph, Graph Convolutional Network

# DGIG-Net: Dynamic Graph-In-Graph Networks for Few-Shot Human-Object Interaction

Xiyao Liu, *Student Member, IEEE*, Zhong Ji, *Member, IEEE*, Yanwei Pang, *Senior Member, IEEE*,  
Jungong Han, *Senior Member, IEEE*, and Xuelong Li, *Fellow, IEEE*

**Abstract**—Few-Shot Learning (FSL) for Human-Object Interaction (HOI) aims at recognizing various relationships between human actions and surrounding objects only from few samples. It is a challenging vision task, in which the diversity and interactivity of human actions result in great difficulty to learn an adaptive classifier to catch ambiguous inter-class information. Therefore, traditional FSL methods usually perform unsatisfactorily in complex HOI scenes. To this end, we propose Dynamic Graph-In-Graph Networks (DGIG-Net), a novel graph prototypes framework to learn a dynamic metric space by embedding a visual sub-graph to a task-oriented cross-modal graph for few-shot HOI. Specifically, we first build a knowledge reconstruction graph to learn latent representations for HOI categories by reconstructing the relationship among visual features, which generates visual representations under the category distribution of every task. Then, a dynamic relation graph integrates both reconstructible visual nodes and dynamic task-oriented semantic information to explore a graph metric space for HOI class prototypes, which applies the discriminative information from the similarities among actions or objects. We validate DGIG-Net on multiple benchmark datasets, on which it largely outperforms existing few-shot learning approaches and achieves state-of-the-art results.

**Index Terms**—Few-Shot Learning, Human-Object Interaction, Meta-Learning, Dynamic Graph, Graph Convolutional Network.

## I. INTRODUCTION

UNDERSTANDING human actions and activities from vision information is a long-standing research for building an intelligent system [1], [2]. One important direction is Human-Object Interaction (HOI), which aims at recognizing various relationships between human actions and surrounding objects. However, the development of the HOI study strikes a bottleneck, in which current techniques are difficult to address the imbalanced data distribution in HOI. Recently, Few-Shot Learning (FSL) provides HOI a novel solution due to its potential to alleviate the low-data challenge.

Few-Shot Learning (FSL) for Human-Object Interaction (HOI) is proposed to recognize novel HOI categories effectively with a limited number of labeled examples [3]. It has

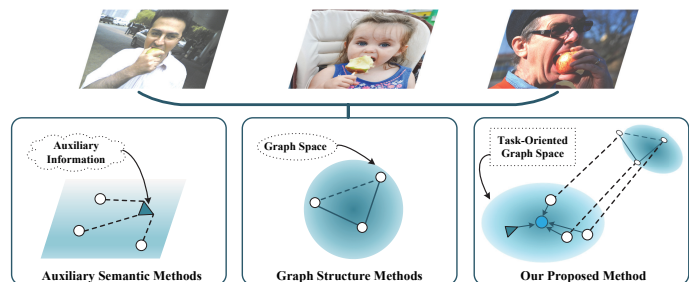


Fig. 1. The difference between our model and other methods. Auxiliary semantic methods introduce cross-modal information to help recognize new visual objects. Graph structure methods build a topological structure to improve the generalization of instance representations. Our proposed method designs a graph-in-graph structure to embed a visual sub-graph to a dynamic graph metric space guided by task-oriented semantic knowledge.

the potential to address: (1) The recognition of tail part in HOI distribution. HOI data is a natural long-tail distribution, where the instance imbalance among categories suffers from over-fitting [4]. FSL methods learn a network that maps an unlabeled example (query sample) to its label from the small labeled support set [5], which imitates the capability of humans to identify objects with very little direct supervision. (2) The combinatorial explosion problem in HOI. Multiple labels in HOI cause the number of classes to increase exponentially, which results in difficulty to solve large scale practical problems. FSL methods transfer the knowledge from existing HOI models to recognize novel visual concepts instead of training a new model from scratch. Although FSL methods provide the HOI scene with a promising direction, the HOI scene brings new challenges to the existing FSL methods.

The purpose of HOI emphasizes the relationships between objects and people, which are quite diverse and interactive. The same action with different objects is classified as different HOI categories, which results in the difference in inter-class is ambiguous. For example, “Eat-Apple” and “Eat-Banana” are different classes in HOI, but they are similar in visual representations. This is a big challenge for current FSL methods, which is still in its infancy and only implements in a simple and single scene, such as miniImageNet [5]. It is difficult to learn an adaptive learner for complex HOI scenes, which results in the unsatisfactory performance of FSL methods. An effective approach to improve the few-shot performance is to learn more representative and discriminative visual features, which provides sufficient evidence to perform classification with few samples.

Manuscript received xxx xx, 2020; revised xxx xx, 2020.

This work was supported by the National Natural Science Foundation of China (NSFC) under Grants 61771329 and 61632018, and the Natural Science Foundation of Tianjin under Grant 19JCYBJC16000.

X. Liu, Z. Ji\*(corresponding author), and Y. Pang are with the School of Electrical and Information Engineering, Tianjin University, and Tianjin Key Laboratory of Brain-inspired Intelligence Technology, Tianjin 300072, China (e-mails: xiyao.liu@tju.edu.cn; jizhong@tju.edu.cn; pyw@tju.edu.cn).

J. Han is with the Data Science group, University of Warwick, Coventry CV4 7AL, UK (e-mail: jungong.han@warwick.ac.uk).

X. L. is with the School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an 710072, China (e-mail: xuelong\_li@ieee.org).

For improving the representativeness and discriminability of instances, recent FSL studies develop two types of approaches. One is introducing auxiliary semantic modalities such as label embeddings [3], [6], attribute annotations [7], and description text [8]. These approaches are inspired by that language explanations help infants to recognize new visual objects, thus providing a strong information source for data scarcity issue [9]. The other is designing a graph-based method owing to its advantage on the effective representation of the graph-structured data. Graph structure models in the few-shot learning have achieved promising results by dealing with complex relationships and inter-dependency among instances [10], [11]. Motivated by the above observations, we design a graph-in-graph structure to embed a visual sub-graph to a dynamic graph metric space guided by task-oriented semantic knowledge, as shown in Fig. 1.

To this end, we propose a Dynamic Graph-In-Graph Networks (DGIG-Net) for few-shot HOI by applying the dynamic and discriminative information from the similarities among actions or objects in this paper, as shown in Fig. 2. It includes a Knowledge Reconstruction Module (KR-Module) and a Dynamic Relation Module (DR-Module), respectively. The KR-Module is designed to reconstruct the relationship among visual features to learn latent representations for HOI categories. Specifically, the encoder exploits both graph structure and node features with a Graph Convolutional Network (GCN), and the decoder reconstructs the topological graph information and manipulates the latent graph representation. The DR-Module implements a graph metric space with dynamic task-oriented semantic information to obtain HOI class prototypes. It applies a cross-modal graph structure to encode two important types of knowledge: (1) The semantic guidance by action and object labels, dynamically defined by the label information from Word2Vector [12]; (2) The visual features obtained by the KR-Module.

It is worthwhile to highlight several aspects of the proposed approach here:

- We implement a novel graph prototypes framework, Dynamic Graph-In-Graph Networks (DGIG-Net) by embedding a visual sub-graph to a dynamic graph metric space. In this way, it obtains HOI class prototypes instead of linear prototypes method, which improves the representativeness and discriminability of the prototype features.
- We design the Knowledge Reconstruction Module (KR-Module) to encode both graph structure and node features with a Graph Convolutional Network (GCN), and reconstruct the topological graph information, which manipulates the latent graph representation with a decoder. The KR-Module reconstructs the relationship among visual features to learn latent representations for HOI categories.
- We develop a Dynamic Relation Module (DR-Module) that applies a cross-modal graph structure to encode dynamic semantic guidance by action and object labels and the visual features obtained by the KR-Module. The DR-Module implements a graph metric space with dynamic task-oriented semantic information to obtain HOI class prototypes.
- Extensive experiments on 2 HOI benchmark datasets with

2 split strategies, i.e., HICO-NN, TUHOI-NN, HICO-NF, and TUHOI-NF, demonstrate the effectiveness of our method. For example, our DGIG-Net improves accuracy by 3.7% in terms of 5-way 1-shot and 2.2% in terms of 5-way 5-shot on HICO-NF, and 5.4% and 2.8% on TUHOI-NF against the state-of-the-art methods, respectively. For the cross-domain few-shot HOI task, it also outperforms state-of-the-art methods.

The remaining sections of the paper are organized as follows. Section II reviews the related work. Section III introduces our proposed DGIG-Net in detail. Section IV presents the experiments and analyses, followed by the conclusion in Section V.

## II. RELATED WORK

Our work is related to three active areas in machine learning: Few-Shot Learning, Human-Object Interaction Recognition, and Graph Convolutional Network.

**Few-Shot Learning.** It is designed to train models for classification from only a handful of samples. There are three main types of methods to tackle the few-shot task: metric-based, optimization-based, and generation-based approaches.

Methods in [5], [13], [14], [15], [16] aim at building metric-based networks by measuring the distance to realize few-shot learning. For example, Matching Networks [5] apply a Recurrent Neural Network (RNN) to accumulate task information in the embedding space of training samples to predict classes for testing samples. Remarkably, it defines the episode training strategy, which is widely applied by the following studies. Prototypical Networks [13] learn a linear prototype space for classes and classify the query image into the nearest class prototypes. Relation Networks [14] utilize neural networks to measure the possibility of two images belong to the same class, which replaces the traditional artificial defining distance measurement method. DN4 [15] designs a local descriptor to learn the exchangeability of visual patterns across the images in the same class and complete image-to-class measurements. TPN [16] proposes to learn a graph construction module to propagate labels from labeled instances to unlabeled test instances.

Optimization-based FSL methods [17], [18] propose to learn a good initialization by adjusting the optimization algorithm and effectively obtain model parameters that can be learned with a few examples. For example, MAML [17] designs a model-agnostic method based on learning easily adaptable model parameters through gradient descent. Based on the idea, many methods extend this work such as Reptile [18] and LEO [19]. For example, Reptile [18] directly implements Stochastic Gradient Descent (SGD) in training instead of computing twice gradients, which requires less computation than MAML.

Moreover, some work addresses the few-shot problem with a data-driven solution, called generation-based methods. These methods create novel samples to augment the training set and improve the performance of current few-shot algorithms. Wang *et al.* [20] trained a Generative Adversarial Network (GAN) to generate new instances, which achieves up to a 6% boost in classification accuracy when only given a single



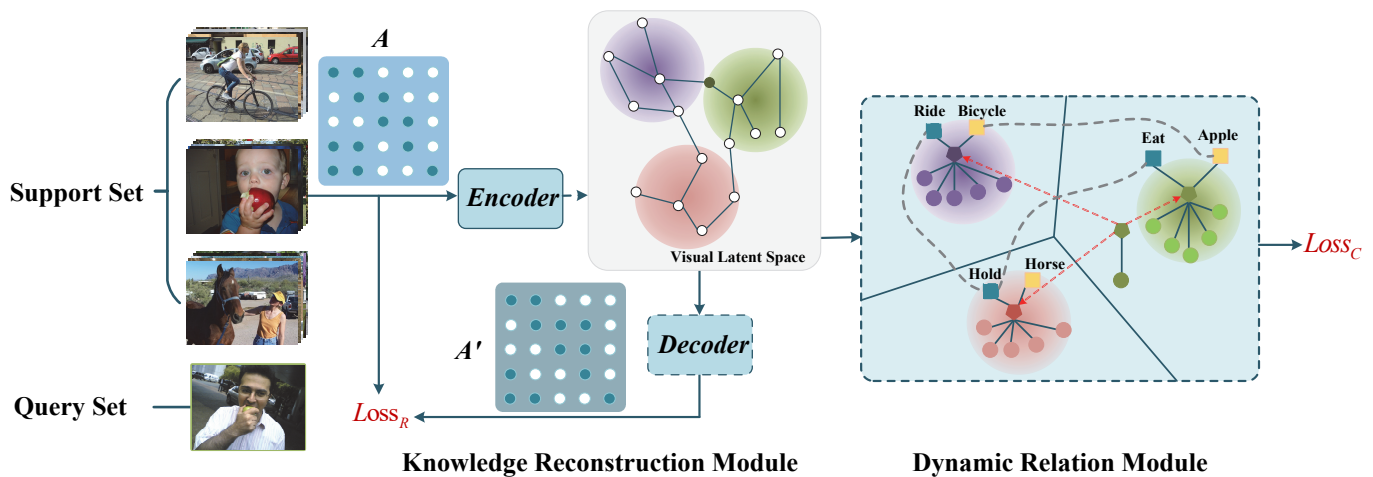


Fig. 2. The proposed DGIG-Net architecture, taking the 3-way 5-shot task for example. Given an adjacency matrix  $A$ , the Knowledge Reconstruction Module learns a latent representation through a graph encoder and then manipulates it input the Dynamic Relation Module, which applies a cross-modal graph structure to obtain HOI class prototypes with dynamic task-oriented semantic information.

training example. Zhang *et al.* [21] replaced backgrounds and foregrounds guided by saliency maps to generate new support samples, which is more effective and less costly than GAN.

Few-shot HOI is a challenging vision task due to the fact that it is difficult to learn an adaptive classifier for complex HOI scenes. For improving the representativeness and discriminability of instances in HOI, we propose a Dynamic Graph-In-Graph Networks (DGIG-Net), a novel graph prototypes framework to learn a dynamic graph metric space guided by task-oriented semantic for few-shot HOI.

**Human-Object Interaction Recognition.** As a sub-task of human action recognition, HOI recognition is introduced for alleviating such ambiguities caused by no motion cue in a still image. Since actions involving objects provide a context in spatial and functional relations, human behaviors are recognized effectively [22]. Early researches relied on shape features and movement analysis to recognize HOI [23].

Recently, deep learning technologies have brought promising results on the HOI recognition task. Successively, large scale image datasets [4], [24] were also released. Some work found that human interacts with an object by contacting some parts of the body instead of all the body. Gkioxari *et al.* [25] developed a part-based model to make fine-grained action recognition based on the input of both whole-person and part bounding boxes. Fang *et al.* [26] proposed a new pairwise body-part attention model that can learn to focus on crucial parts and their correlations for HOI recognition. Moreover, Mallya *et al.* [27] attached importance to HOI in Visual Question Answering (VQA), and proposed a deep convolutional network model that fuses features from local and global context to recognize HOI.

Some studies focuses on instance-based HOI detection tasks [28], [29], which utilizes holistic human poses and global context under the help of popular detectors jointly to infer the locations and categories of HOI. Additionally, compositional learning [30] employs an external knowledge graph to recognize unseen interactions, which applies zero-shot learning

to address the data scarcity problem in HOI. For further alleviating the instance imbalance and combinatorial explosion challenges in HOI recognition, SGAP-Net [3] formulates HOI as a few-shot HOI task and learns a semantic-guided metric space to obtain attentive class prototypes for few-shot HOI.

**Graph Convolutional Network.** Traditional neural networks have made great progress in Euclidean data but still perform unsatisfactorily on non-Euclidean domains. Graph neural networks raise attention by dealing with complex relationships and inter-dependency among instances [31]. Several studies have applied different types of graph neural networks in node classification [30], link prediction [32] and graph classification [33]. For example, Kato *et al.* [30] proposed a zero-shot HOI method, which constructs an external knowledge graph and Graph Convolutional Network (GCN) to recognize novel action-object compositions in HOI. Qi *et al.* [32] inferred a parse graph neural network that includes the HOI graph structure represented by an adjacency matrix, and the node labels for HOI detection. Mallea *et al.* [33] proposed a model for graph classification by extracting fixed size tensorial information from each graph in a given set, and employing a Capsule Network to perform classification.

Graph Convolutional Network (GCN), as a type of node-level graph neural network, performs superior in node regression and node classification tasks. GCN was first proposed in [34], which formulates semi-supervised classification as graph node classification. Recent researches have made incremental improvements over GCN [35], [36], [37]. For example, Adaptive Graph Convolutional Network (AGCN) [35] learns generalized and flexible structural relations unspecified by an arbitrary graph structure. It measures two nodes features by a learnable distance function and constructs a residual graph adjacency matrix, which is fed on data without restrictions on graph degree. Dual Graph Convolutional Network (DGCN) [36] proposes a dual GCN architecture embedding semantic information (i.e., global-consistency-based knowledge) with two graph convolutional layers in parallel.

Recently, graph-based methods have been employed in FZL [38], [39], [40]. For example, Gidaris *et al.* [39] designed a GCN-based denoising autoencoder network by taking as input a set of classification weights corrupted with Gaussian noise to reconstruct the target-discriminative classification weights, which regularizes the weight generating meta-model. [40] focuses on modeling clean and noisy data by a graph per class and predicting class relevance of noisy examples. Although GCN has been directly applied in a number of recent FSL methods, the difference of our DGIG-Net lies in that we apply GCN to explore the cross-modal relationship among semantic information, class prototypes, and visual instances, which learns a graph prototypes metric space to obtain HOI prototypes.

### III. DGIG-NET FOR FEW-SHOT HOI RECOGNITION

In this section, we present our proposed Dynamic Graph-In-Graph Networks (DGIG-Net). The architecture of DGIG-Net consists of a Knowledge Reconstruction Module (KR-Module), and a Dynamic Relation Module (DR-Module), as shown in Fig. 2. We first develop a graph auto-encoder, Knowledge Reconstruction Module (KR-Module), which effectively applies both inter-class and intra-class information to learn a latent representation. Based on the representation, a Dynamic Relation Module (DR-Module) is then proposed to integrate category semantic information and visual information towards cross-modal dynamic prototypes. We first introduce the problem formulation and then report our approach in detail.

#### A. Problem Definition

In few-shot classification, there is a meta-train set  $\mathcal{S} = \{x_i, l_i, y_i\}_{i=1}^N$  that consists of  $N$  samples from  $M$  categories, where each  $x_i \in \mathbb{R}^D$  is a  $D$ -dimensional visual feature vector of the  $i$ -th image,  $l_i$  is its label semantic embeddings, and  $y_i$  is one-hot class label. According to the datasets of HOI, the label for an image combines a pair of the action and the object. Every sample in the meta-train set is randomly divided into the support set or the query set. When training in the support set, the semantic vectors of labels are given as  $l_i = \{(n_i, v_i)\}$ , where  $n_i \in \mathbb{R}^V$  is the  $V$ -dimensional text semantic embedding of the noun label,  $v_i \in \mathbb{R}^V$  is the  $V$ -dimensional text semantic embedding of the verb label. There are no semantic labels in the query set.

We follow a  $C$ -way  $K$ -shot episode-based training strategy defined by Matching Networks [5]. Each episode is formed by sampling  $C$  classes and  $K$  labeled samples of each class from  $\mathcal{S}$  to construct a few-shot task, which contains a support set and a query set to simulate the training and testing process. Specifically, in the  $w$ -th episode, the support set can be denoted as  $\mathcal{S}_{support}^w = \{x_i, l_i, y_i\}_{i=1}^{N_s}$  ( $N_s = C \times K$ ), and the query set  $\mathcal{S}_{query}^w = \{x_i, y_i\}_{i=1}^{N_q}$ .

#### B. Knowledge Reconstruction Module

To represent the relationship of visual space and graph structure in a unified framework, we develop a new graph auto-encoder network as a graph encoder. The idea is to learn

the hidden representations of each node by combining support samples of the same categories, to integrate inter-class visual features with the graph structure in the latent representation. The most straightforward strategy to attend the neighbors of a node is to embed its representation from all its neighbors.

Formally, we define our knowledge reconstruction graph as  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{Z})$ .  $\mathcal{G}$  is an undirected graph with  $\mathcal{V}$  as its nodes.  $\mathcal{E}$  presents the links between nodes  $\mathcal{V}$ , and  $\mathcal{Z}$  are the feature vectors for nodes  $\mathcal{V}$ . Specifically, we deploy an adjacency matrix to represent the visual relationship between support and query samples:

$$A_{KR} = \begin{bmatrix} A_{ss} & 0 \\ 0 & A_{qq} \end{bmatrix}, Z_{KR} = [Z_s, Z_q], \quad (1)$$

where  $A_{KR}$  and  $Z_{KR}$  are the adjacency matrix and node features of the knowledge reconstruction graph,  $A_{ss}$  and  $A_{qq}$  are adjacency matrices for support-support nodes and query-query nodes,  $Z_s$  and  $Z_q$  are visual features of the support set and the query set, respectively. This graph includes 2 types of nodes: support nodes and query nodes. The adjacency matrix  $A_{ss}$  is defined as:

$$A_{ij} = \begin{cases} 1, & \text{if } y_i = y_j \text{ and } i \neq j \\ 0, & \text{else} \end{cases} \quad (2)$$

where  $y_i$  is the label of the  $i$ -th support sample.  $A_{qq} = 0$ , since query nodes only have self-connection to update its node features, which will be added after normalization. To better capture the graph structure, the adjacency matrix is normalized as being real symmetric positive semidefinite [34]. The adjacency  $A_{KR}$  is normalized as:

$$\hat{A}_{KR} = D^{-\frac{1}{2}}(A_{KR} + I)D^{-\frac{1}{2}}, \quad (3)$$

where  $D$  is the diagonal node degree matrix for each block,  $I$  is an identity matrix to add self-connection to each node. The structure in our adjacency matrix could be transformed as:

$$\hat{A}_{KR} = \begin{bmatrix} \hat{A}_{ss} & 0 \\ 0 & \hat{A}_{qq} \end{bmatrix}, \quad (4)$$

where  $\hat{A}_{ss}$  and  $\hat{A}_{qq}$  are adjacency matrices for support-support nodes and query-query nodes after graph normalization.

All visual node features  $Z$  are obtained by transforming the node features that they link on the graph in GCN. Formally, a single layer GCN is calculated as:

$$\tilde{Z} = GCN(Z, A) = \hat{A}Z^TW, \quad (5)$$

where  $\hat{A}$  is the normalized graph adjacency matrix,  $Z$  is the node features,  $W$  is a  $d \times \bar{d}$  weight parameter matrix.  $d$  is the dimensionality of input feature vector for each node and  $\bar{d}$  is the output feature dimensionality. GCN first collects the features of connected nodes with link parameters in  $\hat{A}$ , then transforms features on each node by  $W$  independently. This operation is usually stacked with multi-layer, where non-linear activation functions (i.e., ReLU) are applied.

Since  $\hat{A}$  is a block matrix, it can be further decomposed each GCN layer to each block. This decomposition provides better insights for our model. Specifically, we have:

$$\begin{aligned} \tilde{Z}_s &= \hat{A}_{ss}Z_sW, \\ \tilde{Z}_q &= \hat{A}_{qq}Z_qW, \end{aligned} \quad (6)$$

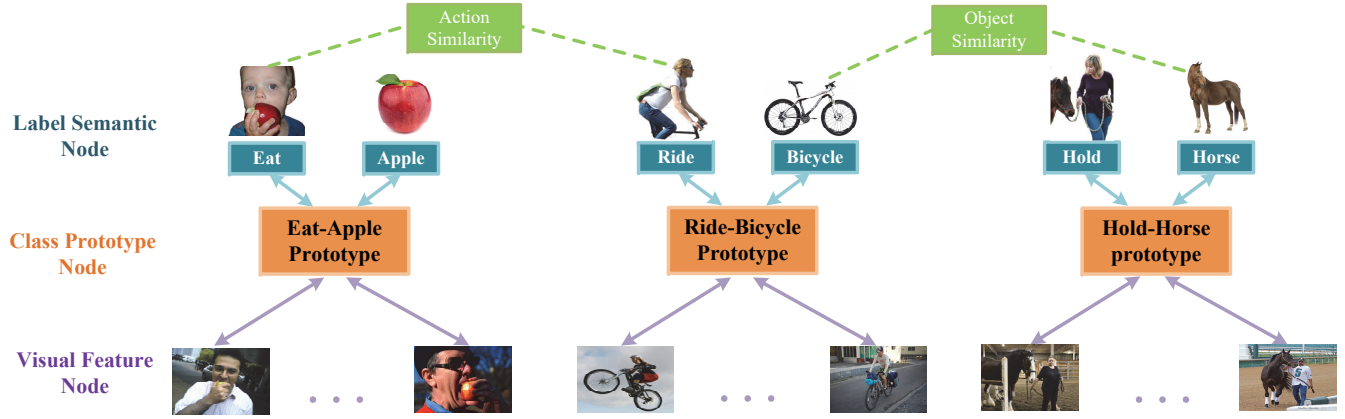


Fig. 3. Illustration of the DR-Module. There are 3 types of nodes in this graph, which are denoted as label semantic nodes, class prototype nodes, and visual feature nodes, respectively. The label semantic nodes provide dynamic relationships from their corresponding task, which are integrated into class prototype nodes by GCN encoder.

where  $\tilde{Z}_s$  and  $\tilde{Z}_q$  are the outputs of our knowledge reconstruction graph that can also be noted by  $Z_v = [\tilde{Z}_s, \tilde{Z}_q]$ . There are various types of decoders, which reconstruct either the graph structure or the attribute value [41]. As our latent embedding already contains both content and structure information, we apply a simple inner product decoder to predict the links between nodes, which would be efficient and flexible:

$$A'_{ij} = \text{Sigmoid}(z_i^T z_j), \quad (7)$$

where  $A'$  is the reconstructed structure matrix of the graph. We minimize the reconstruction error by measuring the difference between  $A$  and  $A'$ :

$$\text{Loss}_R = KL(A' || A) = \sum_{i=1}^n \sum_{j=1}^n A'_{ij} \log \frac{A'_{ij}}{A_{ij}}. \quad (8)$$

### C. Dynamic Relation Module

One of the main challenges for few-shot methods is the dynamic adaption for different tasks. To confront this challenge, we develop a dynamic relation graph algorithm as the solution, which generates a task-oriented graph prototypes metric space.

We design a cross-modal graph structure to encode two important types of knowledge: (1) The semantic guidance by the verb and noun labels, defined by the label information from Word2Vector [12]; (2) The visual features obtained by the KR-Module.

**Graph Construction.** Specifically, we construct the graph as follows.

- (a) There are 3 types of nodes in our graph. These nodes are denoted as label semantic nodes, class prototype nodes and visual feature nodes, where their node features are denoted as  $Z_l$ ,  $Z_p$  and  $Z_v$ , respectively.
- (b) Each class prototype node defines an HOI class prototype. These class prototypes are modeled by a separate set of label semantic nodes  $Z_l$  and visual nodes  $Z_v$  in the graph. These nodes are initialized with all zero feature vectors and will obtain their representations  $Z_p$  via integrating category semantic information and visual information.

- (c) A label semantic (verb or noun) node only connects to a class prototype node. Similarly, each visual (support or query) node only links to a class prototype node.
- (d) WordNet [42] is applied to create noun-noun and verb-verb links, which generates dynamic relationships for the whole graph by the multi-layer GCN encoder.

The graph construction of DR-Module is shown in Fig. 3. This graph is thus captured by its adjacency matrix  $A \in \mathbb{R}^{|V| \times |V|}$  and a feature matrix  $Z \in \mathbb{R}^{d \times |V|}$ . Based on the construction, our graph structure can be naturally decomposed into blocks, given by:

$$A_{DR} = \begin{bmatrix} A_{ll} & A_{lp} & 0 \\ A_{lp}^T & A_{pp} & A_{pv} \\ 0 & A_{pv}^T & A_{vv} \end{bmatrix}, Z_{DR} = [Z_l, Z_p, Z_v], \quad (9)$$

where  $A_{ll}$ ,  $A_{lp}$ ,  $A_{pp}$ ,  $A_{pv}$  and  $A_{vv}$  are adjacency matrices for label-label pairs, label-prototype pairs, prototype-prototype pairs, prototype-visual pairs and visual-visual pairs, respectively.  $Z_l$ ,  $Z_p$  and  $Z_v$  are node features. Especially, a nonlinear network is applied to embed semantic vectors to visual features space to obtain  $Z_l$ . Moreover, we have  $Z_p = 0$  and thus the prototypes nodes need to learn new representations for recognition.  $A_{ll}$  is the adjacency matrix containing the relationship between action label nodes and object label nodes, which is demoted as:

$$A_{ll} = \begin{bmatrix} A_{aa} & A_{ao} \\ A_{ao}^T & A_{oo} \end{bmatrix}, \quad (10)$$

where  $A_{aa}$ ,  $A_{ao}$  and  $A_{oo}$  are adjacency matrices for action-action label pairs, action-object label pairs, and object-object label pairs, respectively. They are defined by the two aspects: (1) The existing action-object task of this episode; (2) The similarity of words calculated by WordNet [42]. Thus,  $A_{ll}$  is dynamically decided by different tasks, which makes  $A_{DR}$  dynamic.

Similarly, the adjacency matrix  $A_{DR}$  needs to be normalized by Eq (3). The structure in our adjacency matrix could be

transformed as:

$$\hat{A}_{DR} = \begin{bmatrix} \hat{A}_{ll} & \hat{A}_{lp} & 0 \\ \hat{A}_{lp}^T & \hat{A}_{pp} & \hat{A}_{pv} \\ 0 & \hat{A}_{pv}^T & \hat{A}_{vv} \end{bmatrix}, \quad (11)$$

where  $\hat{A}_{pp}$ ,  $\hat{A}_{pv}$  and  $\hat{A}_{vv}$  are adjacency matrices for prototype-prototype nodes, prototype-visual nodes and visual-visual nodes after graph normalization.

The prototype node features  $Z_p$  are obtained by transforming the node features that they link on the graph in GCN. Formally, a single layer GCN is calculated by Eq (5).

After the first layer GCN encoder, we have:

$$\begin{aligned} \tilde{Z}_l &= (\hat{A}_{ll}Z_l + \hat{A}_{lp}Z_p)W, \\ \tilde{Z}_p &= (\hat{A}_{lp}^TZ_l + \hat{A}_{pp}Z_p + \hat{A}_{pv}Z_v)W, \\ \tilde{Z}_v &= (\hat{A}_{pv}Z_p + \hat{A}_{vv}Z_v)W, \end{aligned} \quad (12)$$

where  $\tilde{Z}_l$ ,  $\tilde{Z}_p$  and  $\tilde{Z}_v$  are the outputs of our HOI graph that can also be noted by  $\tilde{Z} = [\tilde{Z}_l, \tilde{Z}_p, \tilde{Z}_v]$ . With nonlinear activations and multi-layer GCN, the model will construct a nonlinear transform that considers more nodes for building the HOI class prototypes. We implement the output HOI prototype representations  $\tilde{Z}_p$  for the HOI class prototypes in few-shot HOI.

And we calculate a probability distribution by the distance between a query sample and the class prototypes of support set to accomplish the recognition task:

$$p_\phi(y = c | q \in \mathcal{S}_{query}^w) = \frac{\exp(-d(z_p^q, z_p^c))}{\sum_k \exp(d(z_p^q, z_p^k))}, \quad (13)$$

where  $d(\cdot)$  is the Euclidean distance,  $z_p^c$  is the prototype features of class  $c$  in  $\tilde{Z}_p$  and  $z_p^q$  is the query  $q$  representation after GCN in  $\tilde{Z}_p$ .

Besides, in the training process, we apply a cross-entropy loss to measure the classification error:

$$Loss_C = d(z_p^q, z_p^c) + \log \sum_{N_C} \exp(-d(z_p^q, z_p^c)). \quad (14)$$

Thus, the final loss of the whole DGIG-Net consists of two parts: reconstruction loss and classification loss, which is denoted as:

$$Loss = Loss_C + \lambda Loss_R, \quad (15)$$

where  $\lambda$  is a hyperparameter adjusting the weight of the two modules.

#### IV. EXPERIMENTS

##### A. Experiment Setup

The CNN structure of our model is a pre-trained ResNet-18 [43]. Thus, each input image is represented as a 1,000-dimensional vector. The Adam optimizer is utilized with the initial learning rate is 0.000001. The hyperparameter  $\lambda$  is set to be 0.1. In terms of the regularizer, we set 0.01 for all datasets. Besides, we apply Word2vector [12] to extract the semantic embeddings for the category labels, which are represented as 400-dimensional vectors.

##### B. Datasets

Among the widely available datasets for HOI, we select two popular datasets, namely Humans Interacting with Common Objects (HICO) [4] and Trento Universal Human-Object Interaction (TUHOI) [24]. HICO dataset consists of 42,109 images with 80 objects and 92 actions, which covers almost human daily activities with 377 interactions. For establishing a more natural and realistic dataset, TUHOI collects images first and then defines actions from images instead of some predefined human actions. Thus, TUHOI is a small scale but rich dataset, which contains 9802 images with 95 objects, 66 actions and 194 interactions.

For satisfying the need for our experiments, original datasets should be divided into novel compositions. Following the popular setting of FSL [5], [13], we apply 60/20/20 training/validation/testing repartitions for reorganizing the datasets. And we present 2 split strategies: 1) Novel Noun (NN), and 2) Novel Few instances (NF). Details of both strategies are described below.

1) *Novel Noun*: This is the first split strategy. We follow that ubiquitous similarity exists in the same action interacting with different objects. Moreover, objects could guide a set of behaviors, i.e., apple can be eaten, held and etc. Similar actions with different objects can be transferable knowledge. Thus, we divide objects as different tasks in our work. Specifically, we divide all noun labels into the meta-train set, the meta-val set, and the meta-test set that are disjoint in nouns. For example, the object “apple” only appears in the meta-train set, corresponding similar object “banana” is divided into the meta-test set, as shown in Fig. 4.

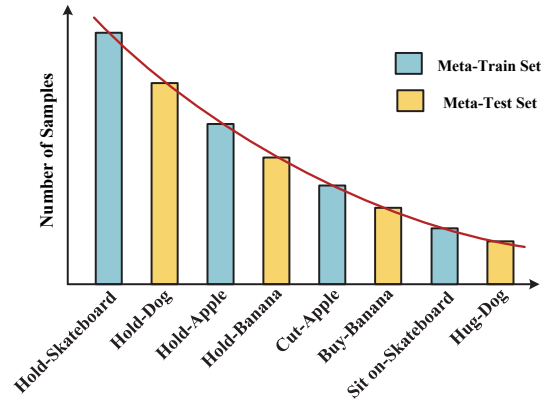


Fig. 4. NN-datasets divide objects as different tasks. The long-tail distribution with disjoint nouns is preserved in each set (taking HICO-NN as an example).

**HICO-NN.** The modified dataset HICO divided with the novel noun strategy is called HICO-NN. We divide HICO-NN into a meta-train set with 45 nouns and 24,067 images, a meta-test set with 20 nouns and 9,146 images, and a meta-validation set with 15 nouns and 8,896 images, which are disjoint in noun labels.

**TUHOI-NN.** TUHOI-NN is reorganized similarly to that of the HICO-NN dataset. There are 50 nouns and 4871 images in the meta-train set, 20 nouns and 2361 images in the meta-val set, and 25 nouns and 2570 images in the meta-test set of TUHOI-NN. More details are listed in Table I.

TABLE I  
Settings of NN-datasets.

Dataset	Item	Meta-train Set	Meta-Val Set	Meta-Test Set	Total
HICO-NN	Action	69	42	47	-
	Object	45	15	20	80
	Interaction	212	73	92	377
	Image	24,067	8,896	9,146	42,109
TUHOI-NN	Action	46	24	24	-
	Object	50	20	25	95
	Interaction	98	44	52	194
	Image	4,871	2,361	2,570	9,802

TABLE II  
Settings of NF-datasets.

Dataset	Item	Meta-train Set	Meta-Val Set	Meta-Test Set	Total
HICO-NF	Action	50	54	49	-
	Object	65	61	59	-
	Interaction	173	102	102	377
	Image	38,147	1,987	1,975	42,109
TUHOI-NF	Action	30	32	28	-
	Object	58	53	53	-
	Interaction	69	62	63	194
	Image	7,294	1,277	1,231	9,802

2) *Novel Few Instances*: This is the second split strategy. Since the data shows a long tail distribution and our purpose is to recognize the unusual categories, we select the categories that have over 50 samples as the meta-train, and others are employed in the meta-test set and the meta-validation set, as shown in Fig. 5.

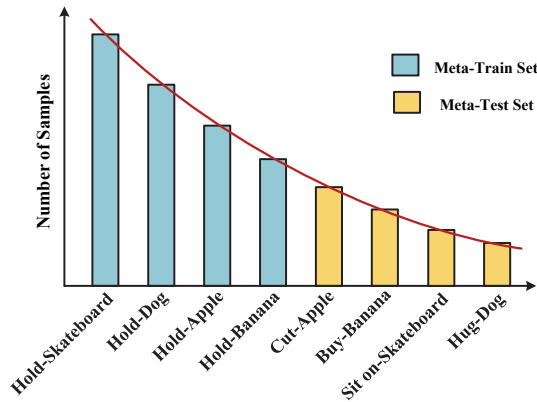


Fig. 5. NF-datasets divide categories with more samples as the meta-train and others are employed in the meta-test set and the meta-validation set. The head part of the long-tail distribution appears in the meta-train and the meta-test set presents the tail part (taking HICO-NF as an example).

**HICO-NF.** We call the modified dataset HICO for novel few instances as HICO-NF. There are 173 interactions and 38,147 images in the meta-train set, 102 interactions and 1,987 images in the meta-validation set, and 102 interactions and 1,975 images in the meta-test set. The categories in the meta-validation and meta-test set contain 49 samples at most and 6 samples at least.

**TUHOI-NF.** Similar to that of HICO-NF dataset, we reorganize TUHOI to be TUHOI-NF. We divide it into a meta-train set with 69 interactions and 7,294 images, a meta-validation set with 62 interactions and 1,277 images, and a meta-test set with 63 interactions and 1,231 images. More details are listed in Table II.

### C. Comparison with State-of-the-Art Methods

We compared a total of 8 few-shot approaches with our model in our experiments. These few-shot algorithms include metric-based and optimization-based methods as follows:

#### Metric-Based Methods:

Matching Networks [5] apply a Recurrent Neural Network (RNN) to accumulate task information in the embedding space.

Prototypical Networks [13] train a CNN to embed task examples to a metric space and perform nearest neighbor classification with the class prototypes.

Relation Networks [14] introduce a learnable metric network to compare the similarity of different samples.

DN4 [15] designs a local descriptor to learn the exchangeability of visual patterns across the images in the same class, and complete based image-to-class measures.

TPN [16] proposes to learn a graph construction module to propagate labels from labeled instances to unlabeled test instances.

SGAP-Net [3] is the first approach designed for few-shot HOI, which learns a semantic-guided metric space to obtain attentive class prototypes.

#### Optimization-Based Methods:

MAML [17] provides a parameter-optimization method for an arbitrary learner model that can be quickly adapted to a particular task.

Reptile [18] generalizes first-order MAML and ignores second-order derivatives, which requires less computation and memory than MAML.

These approaches all utilize ResNet-18 [43] as the embedding networks. It is computed by averaging 10 times over 600 randomly generated episodes as few-shot HOI recognition accuracy.

1) *Comparison on NN Split Strategy*: Table III describes the classification performance of DGIG-Net and eight competitors on HICO-NN and TUHOI-NN. We observe that our approach beats the state-of-the-art in terms of both 5-way 5-shot and 5-way 1-shot tasks on both datasets. Specifically, compared with the second-best method SGAP-Net, the accuracy improvement on HICO-NN in terms of 5-way 1-shot increases from 38.16% to 39.13%, and in terms of 5-way 5-shot from 58.39% to 59.06%. On TUHOI-NN datasets, the proposed DGIG-Net also gains improvements from 37.27% to 38.77% in terms of 5-way 1-shot and from 57.05% to 58.07% in terms of 5-way 5-shot, which outperform the state-of-the-art approaches at least in 1.5% and 1.0%. Moreover, compared with the metric-based approaches, our DGIG-Net achieves obvious improvements on both datasets, which demonstrates graph prototypes capture a more discriminative metric space than the others. Compared with the optimization-based methods, our DGIG-Net has more significant improvements, which



TABLE III  
Few-shot classification accuracy of DGIG-Net on HICO-NN and TUHOI-NN with  $\pm 95\%$  confidence intervals.

Method	Type	HICO-NN		TUHOI-NN	
		5-way 1-shot	5-way 5-shot	5-way 1-shot	5-way 5-shot
Matching Networks [5]	Metric	32.14 $\pm$ 1.62%	44.87 $\pm$ 1.74%	32.48 $\pm$ 1.58%	40.04 $\pm$ 1.70%
Prototypical Networks [13]	Metric	32.56 $\pm$ 1.59%	42.49 $\pm$ 1.75%	31.12 $\pm$ 1.55%	39.26 $\pm$ 1.71%
Relation Networks [14]	Metric	33.20 $\pm$ 1.68%	46.15 $\pm$ 1.81%	33.50 $\pm$ 1.68%	41.15 $\pm$ 1.75%
DN4 [15]	Metric	33.07 $\pm$ 1.43%	46.19 $\pm$ 1.74%	32.49 $\pm$ 1.43%	41.75 $\pm$ 1.77%
TPN [16]	Metric	33.40 $\pm$ 1.55%	46.33 $\pm$ 1.86%	32.95 $\pm$ 1.59%	41.73 $\pm$ 1.79%
SGAP-Net [3]	Metric	38.16 $\pm$ 1.65%	58.39 $\pm$ 1.82%	37.27 $\pm$ 1.61%	57.05 $\pm$ 1.73%
MAML [17]	Optimization	33.87 $\pm$ 1.74%	47.25 $\pm$ 1.84%	33.78 $\pm$ 1.64%	43.67 $\pm$ 1.79%
Reptile [18]	Optimization	33.26 $\pm$ 1.77%	46.56 $\pm$ 1.85%	32.39 $\pm$ 1.81%	41.65 $\pm$ 1.93%
<b>DGIG-Net (Ours)</b>	Metric	<b>39.13 <math>\pm</math> 1.68%</b>	<b>59.06 <math>\pm</math> 1.89%</b>	<b>38.77 <math>\pm</math> 1.49%</b>	<b>58.07 <math>\pm</math> 1.89%</b>

TABLE IV  
Few-shot classification accuracy of DGIG-Net on HICO-NF and TUHOI-NF with  $\pm 95\%$  confidence intervals.

Method	Type	HICO-NF		TUHOI-NF	
		5-way 1-shot	5-way 5-shot	5-way 1-shot	5-way 5-shot
Matching Networks [5]	Metric	34.95 $\pm$ 1.64%	46.01 $\pm$ 1.70%	33.00 $\pm$ 1.63%	41.59 $\pm$ 1.72%
Prototypical Networks [13]	Metric	34.88 $\pm$ 1.60%	45.47 $\pm$ 1.76%	32.04 $\pm$ 1.57%	41.27 $\pm$ 1.75%
Relation Networks [14]	Metric	36.62 $\pm$ 1.68%	48.01 $\pm$ 1.89%	36.05 $\pm$ 1.67%	42.35 $\pm$ 1.73%
DN4 [15]	Metric	35.21 $\pm$ 1.43%	47.36 $\pm$ 1.74%	35.47 $\pm$ 1.43%	44.72 $\pm$ 1.77%
TPN [16]	Metric	36.24 $\pm$ 1.37%	49.35 $\pm$ 1.66%	37.67 $\pm$ 1.56%	45.32 $\pm$ 1.67%
SGAP-Net [3]	Metric	43.37 $\pm$ 1.66%	70.78 $\pm$ 1.81%	41.12 $\pm$ 1.64%	69.47 $\pm$ 1.74%
MAML [17]	Optimization	38.86 $\pm$ 1.69%	53.32 $\pm$ 1.90%	36.45 $\pm$ 1.84%	48.48 $\pm$ 1.90%
Reptile [18]	Optimization	38.49 $\pm$ 1.69%	52.98 $\pm$ 1.78%	37.26 $\pm$ 1.83%	49.29 $\pm$ 1.88%
<b>DGIG-Net (Ours)</b>	Metric	<b>47.08 <math>\pm</math> 1.44%</b>	<b>73.06 <math>\pm</math> 1.62%</b>	<b>46.54 <math>\pm</math> 1.49%</b>	<b>72.36 <math>\pm</math> 1.62%</b>

indicates that our model with task-oriented dynamic relation is more transferable than learning the optimization strategy.

We also observe that the results of all methods on TUHOI-NN are lower than those on HICO-NN. We suppose the reason lies in the original distribution of data: the average samples of TUHOI-NN are much less than those of HICO-NN. Therefore, the few-shot HOI task is more difficult on TUHOI than that on HICO. Remarkably, DGIG-Net achieves 38.77% in terms of 5-way 1-shot on TUHOI-NN datasets, which significantly outperforms the second-best performance by 1.5%. It demonstrates that the dynamic relation graph structure of DGIG-Net has the superior ability on the difficult dataset with fewer samples. Moreover, our work applies task-oriented semantic guidance to capture class discriminative information, which learns a dynamic graph prototypes metric space in few-shot HOI.

2) *Comparison on NF Split Strategy*: The results on the NF datasets are summarized in Table IV. Our DGIG-Net achieves the accuracies of 47.08% on 5-way 1-shot and 73.06% on 5-way 5-shot on HICO-NF, which outperform the state-of-the-art approaches at least in 3.7% and 2.2%. Similar results are also observed on TUHOI-NF. It can also be observed that the performance on HICO-NF and TUHOI-NF is better than those on HICO-NN and TUHOI-NN. We consider the reason is as follows. From the perspective of data structure, the NN split strategy makes the meta-train set and the meta-test set both follow the similar long-tail distribution, which brings difficulty to transfer knowledge among imbalanced class distributions. In contrast, the NF split strategy divides the meta-train set as a head distribution, and the meta-test set as a tail distribution. It

provides much more knowledge from instances. Moreover, the objects and actions in the test set in the NF split strategy may also appear in the meta-train set separately, as the purpose is to recognize unseen combinations.

#### D. Ablation Studies

TABLE V  
Ablation studies of DGIG-Net on HICO-NN.

Methods	5-way 1-shot	5-way 5-shot
PN	32.56 $\pm$ 1.59%	42.49 $\pm$ 1.75%
Graph PN	35.69 $\pm$ 1.56%	52.34 $\pm$ 1.81%
Graph PN + Actions	36.49 $\pm$ 1.67%	54.38 $\pm$ 1.79%
Graph PN + Actions + R_A	36.86 $\pm$ 1.63%	54.89 $\pm$ 1.81%
Graph PN + Objects	36.62 $\pm$ 1.73%	54.77 $\pm$ 1.74%
Graph PN + Objects + R_O	37.23 $\pm$ 1.65%	55.13 $\pm$ 1.86%
Graph PN + Actions + Objects	38.24 $\pm$ 1.60%	56.39 $\pm$ 1.83%
DGIG-Net w/o KR-Module	39.02 $\pm$ 1.66%	58.38 $\pm$ 1.84%
<b>DGIG-Net</b>	<b>39.13 <math>\pm</math> 1.68%</b>	<b>59.06 <math>\pm</math> 1.89%</b>

We conduct ablation studies to evaluate the impacts of each component in our DGIG-Net in Table VI. We first consider the following variants:

**PN** is Prototypical Networks [13], which is employed as the baseline for DGIG-Net.

**Graph PN** applies graph prototypes metric space instead of linear prototypes in PN.

**Graph PN + Actions** adds action label semantic in Graph PN.

TABLE VI

Cross-domain few-shot classification accuracy of DGIG-Net on HICO-NN→TUHOI-NN and TUHOI-NN→HICO-NN with  $\pm 95\%$  confidence intervals.

Method	Type	HICO-NN→TUHOI-NN		TUHOI-NN→HICO-NN	
		5-way 1-shot	5-way 5-shot	5-way 1-shot	5-way 5-shot
Matching Networks [5]	Metric	29.13 $\pm$ 1.48%	38.83 $\pm$ 1.78%	31.53 $\pm$ 1.62%	38.13 $\pm$ 1.67%
Prototypical Networks [13]	Metric	28.20 $\pm$ 1.55%	38.93 $\pm$ 1.70%	30.13 $\pm$ 1.61%	38.97 $\pm$ 1.75%
Relation Networks [14]	Metric	28.60 $\pm$ 1.69%	35.27 $\pm$ 1.72%	30.97 $\pm$ 1.63%	36.10 $\pm$ 1.73%
DN4 [15]	Metric	35.05 $\pm$ 0.93%	46.82 $\pm$ 1.79%	28.51 $\pm$ 0.83%	37.81 $\pm$ 1.70%
TPN [16]	Metric	35.47 $\pm$ 1.55%	43.35 $\pm$ 1.86%	31.17 $\pm$ 1.48%	38.26 $\pm$ 1.78%
SGAP-Net [3]	Metric	37.96 $\pm$ 1.68%	56.45 $\pm$ 1.77%	36.89 $\pm$ 1.74%	56.35 $\pm$ 1.79%
MAML [17]	Optimization	36.43 $\pm$ 1.69%	47.23 $\pm$ 1.87%	32.87 $\pm$ 1.65%	39.30 $\pm$ 1.80%
Reptile [18]	Optimization	37.54 $\pm$ 1.73%	46.39 $\pm$ 1.85%	33.21 $\pm$ 1.51%	39.67 $\pm$ 1.79%
<b>DGIG-Net (Ours)</b>	Metric	<b>39.09 <math>\pm</math> 1.52%</b>	<b>58.17 <math>\pm</math> 1.87%</b>	<b>38.53 <math>\pm</math> 1.68%</b>	<b>56.86 <math>\pm</math> 1.76%</b>

TABLE VII

Cross-domain few-shot classification accuracy of DGIG-Net on HICO-NF→TUHOI-NF and TUHOI-NF→HICO-NF with  $\pm 95\%$  confidence intervals.

Method	Type	HICO-NF→TUHOI-NF		TUHOI-NF→HICO-NF	
		5-way 1-shot	5-way 5-shot	5-way 1-shot	5-way 5-shot
Matching Networks [5]	Metric	36.32 $\pm$ 1.64%	53.55 $\pm$ 1.72%	33.13 $\pm$ 1.57%	47.26 $\pm$ 1.70%
Prototypical Networks [13]	Metric	35.05 $\pm$ 1.59%	51.26 $\pm$ 1.74%	30.30 $\pm$ 1.56%	47.70 $\pm$ 1.68%
Relation Networks [14]	Metric	41.21 $\pm$ 1.73%	53.78 $\pm$ 1.85%	36.85 $\pm$ 1.61%	47.71 $\pm$ 1.77%
DN4 [15]	Metric	38.05 $\pm$ 1.69%	55.67 $\pm$ 1.98%	32.56 $\pm$ 1.45%	49.81 $\pm$ 1.73%
TPN [16]	Metric	39.47 $\pm$ 1.55%	53.35 $\pm$ 1.86%	36.57 $\pm$ 1.73%	49.99 $\pm$ 1.98%
SGAP-Net [3]	Metric	42.17 $\pm$ 1.63%	70.49 $\pm$ 1.92%	43.29 $\pm$ 1.59%	69.68 $\pm$ 1.87%
MAML [17]	Optimization	41.28 $\pm$ 1.34%	56.34 $\pm$ 1.56%	37.88 $\pm$ 1.39%	54.29 $\pm$ 1.84%
Reptile [18]	Optimization	41.79 $\pm$ 1.77%	57.39 $\pm$ 1.81%	38.97 $\pm$ 1.52%	54.69 $\pm$ 1.78%
<b>DGIG-Net (Ours)</b>	Metric	<b>46.86 <math>\pm</math> 1.56%</b>	<b>72.69 <math>\pm</math> 1.69%</b>	<b>45.47 <math>\pm</math> 1.65%</b>	<b>71.56 <math>\pm</math> 1.74%</b>

**Graph PN + Actions + R\_A** introduces action label semantic and their dynamic relationship obtained by WordNet [42] in Graph PN.

**Graph PN + Objects** adds object label semantic in Graph PN.

**Graph PN + Objects + R\_O** introduces object label semantic and their dynamic relationship obtained by WordNet [42] in Graph PN.

**Graph PN + Actions + Objects** adds both action and object label semantic in Graph PN.

**DGIG-Net w/o KR-Module** applies both object and action label semantic and their dynamic relationship obtained by WordNet [42] in Graph PN.

Firstly, it is observed that Graph PN improves the results at least 3.1% and 9.8% respectively on 5-way 1-shot and 5-way 5-shot compared with PN, as shown in Table V. It proves that graph prototypes method is more effective than linear prototypes method for the few-shot HOI recognition. We can observe that applying a single type of semantic information, i.e., Graph PN + Actions or Graph PN + Objects, brings at least 0.8% and 2.0% performance gains in terms of both settings respectively. This is a reasonable phenomenon since introducing auxiliary semantic information helps to learn a discriminative metric space. By contrast, Graph PN + Actions + Objects applies both types of semantic information, which achieves the surprising accuracies of 38.24% in terms of 5-way 1-shot and 56.39% in terms of 5-way 5-shot. The role of dynamic relationship is proved by Graph PN + Actions + R and Graph PN + Objects + R. They slightly improve the corresponding baseline performance by 0.4% and 0.6%.

Remarkably, DGIG-Net (w/o KR-Module), just Graph PN + Actions + Objects + R\_A + R\_O, achieves 39.02% in terms of 5-way 1-shot and 58.38% in terms of 5-way 5-shot on HICO-NN datasets, which marginally improves the Graph PN + Actions + Objects by 0.7% and 2.0%. It also can be observed that the KR-Module are respectively capable of bringing 0.1% and 0.6% performance gains on both settings against DGIG-Net (w/o KR-Module). The KR-Module seems no more improvement on 5-way 1-shot due to that intra-class information doesn't work on only 1 sample.

#### E. Cross-Domain Analysis

To prove the transferability of the proposed approaches, we design cross-domain experiments that are conducted between two datasets with the same split strategy. There are 4 types of cross-domain settings: HICO-NN→TUHOI-NN, TUHOI-NN→HICO-NN, HICO-NF→TUHOI-NF, and TUHOI-NF→HICO-NF. For the cross-domain setting A→B, the meta-train set of A is utilized in the training stage, while meta-validation and meta-test of B are utilized for validation and evaluation. This cross-domain setting explores that if meta-learning could implement on data from the different source domain and target domain.

We choose the same comparison algorithm listed in Table III. The results of cross-domain experiments on the NN setting are shown in Table VI. Our DGIG-Net achieves competitive performance, which respectively obtains the accuracies of 39.09% and 58.17% on HICO-NN→TUHOI-NN, 38.53% and 56.86% on TUHOI-NN→HICO-NN on 5-way 1-shot and 5-way 5-shot. DGIG-Net performs superior to the second-best



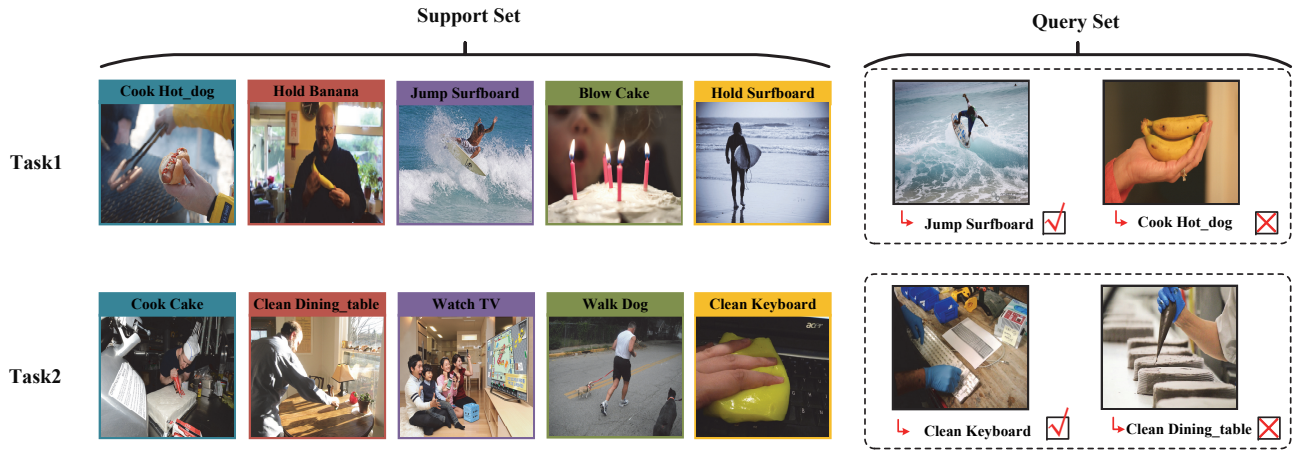


Fig. 6. Qualitative results of few-shot HOI on HICO-NN dataset. The task setting is 5-way 1-shot. (Best viewed in color.)

approach SGAP-Net, which brings about at least 1.0% and 1.6% performance gains on both 5-way 1-shot and 5-way 5-shot tasks on HICO-NN $\rightarrow$ TUHOI-NN. However, our DGIG-Net improves no significant performance in terms of 5-way 5-shot on TUHOI-NN $\rightarrow$ HICO-NN. It can be further observed that the accuracies of HICO-NN on cross-domain settings are a bit lower than those on the single domain experiments. We suppose there are fewer instances in the meta-train set of TUHOI-NN than those of HICO-NN, thus TUHOI-NN can not provide enough transferable knowledge. However, the performance on HICO-NN $\rightarrow$ TUHOI-NN is superior to those on TUHOI-NN, except for the result of DGIG-Net on 5-way 5-shot. It suggests that the source domain with more samples and categories has a higher adaptation ability on the target domain. Our DGIG-Net achieves incremental improvement, which demonstrates its strong adaptable ability.

We also conduct experiments on the NF setting, which are shown in Table VII. Our proposed approaches achieve at least 4.6% and 2.2% gains on HICO-NF $\rightarrow$ TUHOI-NF, and 2.2% and 1.8% gains on TUHOI-NF $\rightarrow$ HICO-NF compared with those of the second-best SGAP-Net [3]. For the comparison between the cross-domain and single domain experiments, the accuracy of TUHOI-NF $\rightarrow$ HICO-NF achieves 71.56% on 5-way 5-shot, which is 1.5% inferior to that of HICO-NF on the single domain experiments. Obviously, it can be observed that the results on HICO-NF are better than those on TUHOI-NF $\rightarrow$ HICO-NF. In Contrast, the performance of DGIG-Net on HICO-NF $\rightarrow$ TUHOI-NF decreases a little, but the other results are better than those of single domain experiments on TUHOI-NF. From our previous analysis, it depends on the data distribution and data scale of different domains. The cross-domain experiments on both NN and NF settings show that the proposed approaches have robust domain adaptation ability against the other approaches.

#### F. Qualitative Results

To further qualitatively verify the effectiveness of our proposed model, we select several representative tasks to show their corresponding few-shot results on HICO-NN. Figure

6 presents the qualitative results on 5-way 1-shot with 2 query samples. It can be observed that our model recognizes HOI categories correctly when appearing the same actions or objects in the task. For example, in Task 1, our model can recognize the query image is “Jump-Surfboard” without the interfere of “Hold-Surfboard”. On the other side, we are trying to explore the reason for misclassification. Concretely, our model performs unsatisfactorily in: 1) Support and query samples present totally different visual angles, such as the two samples of “Hold-Banana” in Task 1. It exists a huge visual bias between local and global views. 2) Objects of different shapes. The “Cake” in Task 2 brings coarse-grained and fine-grained level shapes, which requires more detailed information.

#### V. CONCLUSION

Few-Shot Learning for HOI is a challenging vision task, in which diversity and interactivity of human actions result in great difficulty to learn adaptive class prototypes. In this work, we have proposed a novel graph prototypes framework, namely DGIG-Net, to learn a dynamic graph metric space guided by task-oriented semantic for few-shot HOI. The KR-Module encodes both graph structure and node features with a Graph Convolutional Network, where the decoder reconstructs the topological graph information and manipulates the latent graph representation. The DR-Module implements a graph metric space with dynamic task-oriented semantic information to obtain HOI class prototypes. Extensive experiments on four few-shot HOI datasets, HICO-NN, TUHOI-NN, HICO-NF, and TUHOI-NF, have demonstrated that our proposed approaches are superior to state-of-the-art approaches.

#### REFERENCES

- [1] L. Wu, Y. Wang, X. Li, and J. Gao, “Deep attention-based spatially recursive networks for fine-grained visual recognition,” *IEEE Transactions on Cybernetics*, vol. 49, no. 5, pp. 1791–1802, 2019.
- [2] Z. Zhang, J. Chen, Q. Wu, and L. Shao, “Gii representation-based cross-view gait recognition by discriminative projection with list-wise constraints,” *IEEE Transactions on Cybernetics*, vol. 48, no. 10, pp. 2935–2947, 2018.

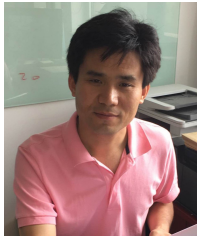
- [3] Z. Ji, X. Liu, Y. Pang, and X. Li, "SGAP-Net: Semantic-guided attentive prototypes network for few-shot human-object interaction recognition," in *AAAI Conference on Artificial Intelligence*, pp. 11085–11092, 2020.
- [4] Y. Chao, Z. Wang, Y. He, J. Wang, and J. Deng, "HICO: A benchmark for recognizing human-object interactions in images," in *International Conference on Computer Vision*, pp. 1017–1025, 2015.
- [5] O. Vinyals, C. Blundell, T. P. Lillicrap, K. Kavukcuoglu, and D. Wierstra, "Matching networks for one shot learning," in *Advances in Neural Information Processing Systems*, pp. 3637–3645, 2016.
- [6] C. Xing, N. Rostamzadeh, B. N. Oreshkin, and P. H. O. Pinheiro, "Adaptive cross-modal few-shot learning," in *Advances in Neural Information Processing Systems*, pp. 4848–4858, 2019.
- [7] P. Tokmakov, Y. Wang, and M. Hebert, "Learning compositional representations for few-shot recognition," in *International Conference on Computer Vision*, pp. 6372–6381, 2019.
- [8] J. Mu, P. Liang, and N. D. Goodman, "Shaping visual representations with language for few-shot classification," *arXiv preprint arXiv:1911.02683*, 2019.
- [9] Y. Yu, Z. Ji, J. Guo, and Z. Zhang, "Zero-shot learning via latent space encoding," *IEEE Transactions on Cybernetics*, vol. 49, no. 10, pp. 3755–3766, 2019.
- [10] H. Li, W. Dong, X. Mei, C. Ma, F. Huang, and B. Hu, "LGM-Net: Learning to generate matching networks for few shot learning," in *International Conference on Machine Learning*, pp. 3825–3834, 2019.
- [11] J. Zhang, M. Zhang, Z. Lu, T. Xiang, and J. Wen, "Adargcn: Adaptive aggregation gcn for few-shot learning," *arXiv preprint arXiv:2002.12641*, 2020.
- [12] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in Neural Information Processing Systems*, pp. 3111–3119, 2013.
- [13] J. Snell, K. Swersky, and R. S. Zemel, "Prototypical networks for few-shot learning," in *Advances in Neural Information Processing Systems*, pp. 4077–4087, 2017.
- [14] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. S. Torr, and T. M. Hospedales, "Learning to compare: Relation network for few-shot learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1199–1208, 2018.
- [15] W. Li, L. Wang, J. Xu, J. Huo, Y. Gao, and J. Luo, "Revisiting local descriptor based image-to-class measure for few-shot learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7260–7268, 2019.
- [16] Y. Liu, J. Lee, M. Park, S. Kim, E. Yang, S. J. Hwang, and Y. Yang, "Learning to propagate labels: Transductive propagation network for few-shot learning," in *International Conference on Learning Representations*, pp. 1–14, 2019.
- [17] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International Conference on Machine Learning*, pp. 1126–1135, 2017.
- [18] A. Nichol, J. Achiam, and J. Schulman, "On first-order meta-learning algorithms," *arXiv preprint arXiv:1803.02999*, 2018.
- [19] A. A. Rusu, D. Rao, J. Sygnowski, O. Vinyals, R. Pascanu, S. Osindero, and R. Hadsell, "Meta-learning with latent embedding optimization," in *International Conference on Learning Representations*, pp. 1–13, 2019.
- [20] Y. Wang, R. B. Girshick, M. Hebert, and B. Hariharan, "Low-shot learning from imaginary data," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7278–7286, 2018.
- [21] H. Zhang, J. Zhang, and P. Koniusz, "Few-shot learning via saliency-guided hallucination of samples," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2770–2779, 2019.
- [22] P. Bach, G. Knoblich, T. C. Gunter, A. D. Friederici, and W. Prinz, "Action comprehension: deriving spatial and functional relations," *Journal of Experimental Psychology: Human Perception and Performance*, vol. 31, no. 3, pp. 465–479, 2005.
- [23] A. Gupta and L. S. Davis, "Objects in action: An approach for combining action understanding and object perception," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8, 2007.
- [24] D. Le, J. R. R. Uijlings, and R. Bernardi, "TUHOI: Trento universal human object interaction dataset," in *Workshop on Vision and Language*, pp. 17–24, 2014.
- [25] G. Gkioxari, R. Girshick, and J. Malik, "Actions and attributes from wholes and parts," in *International Conference on Computer Vision*, pp. 2470–2478, 2015.
- [26] H. Fang, J. Cao, Y. Tai, and C. Lu, "Pairwise body-part attention for recognizing human-object interactions," in *European Conference on Computer Vision*, pp. 52–68, 2018.
- [27] A. Mallya and S. Lazebnik, "Learning models for actions and person-object interactions with transfer to question answering," in *European Conference on Computer Vision*, pp. 414–428, 2016.
- [28] G. Gkioxari, R. B. Girshick, P. Dollar, and K. He, "Detecting and recognizing human-object interactions," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8359–8367, 2018.
- [29] Y. Li, S. Zhou, X. Huang, L. Xu, Z. Ma, H. Fang, Y. Wang, and C. Lu, "Transferable interactiveness knowledge for human-object interaction detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3585–3594, 2019.
- [30] K. Kato, Y. Li, and A. Gupta, "Compositional learning for human-object interaction," in *European Conference on Computer Vision*, pp. 234–251, 2018.
- [31] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, "A comprehensive survey on graph neural networks," *arXiv preprint arXiv:1901.00596*, 2019.
- [32] S. Qi, W. Wang, B. Jia, J. Shen, and S. C. Zhu, "Learning human-object interactions by graph parsing neural networks," in *European Conference on Computer Vision*, pp. 407–423, 2018.
- [33] M. D. G. Mallea, P. Meltzer, and P. J. Bentley, "Capsule neural networks for graph classification using explicit tensorial graph representations," *arXiv preprint arXiv:1902.08399*, 2019.
- [34] T. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, pp. 1–14, 2017.
- [35] R. Li, S. Wang, F. Zhu, and J. Huang, "Adaptive graph convolutional neural networks," in *AAAI Conference on Artificial Intelligence*, pp. 3546–3553, 2018.
- [36] C. Zhuang and Q. Ma, "Dual graph convolutional networks for graph-based semi-supervised classification," in *The Web Conference*, pp. 499–508, 2018.
- [37] X. Zhou, F. Shen, L. Liu, W. Liu, L. Nie, Y. Yang, and H. T. Shen, "Graph convolutional network hashing," *IEEE Transactions on Cybernetics*, vol. 50, no. 4, pp. 1460–1472, 2020.
- [38] V. G. Satorras and J. B. Estrach, "Few-shot learning with graph neural networks," pp. 1–13, 2018.
- [39] S. Gidaris and N. Komodakis, "Generating classification weights with gnn denoising autoencoders for few-shot learning," in *International Conference on Learning Representations*, pp. 21–30, 2019.
- [40] A. Iscen, G. Tolias, Y. Avrithis, O. Chum, and C. Schmid, "Graph convolutional networks for learning with few clean and many noisy labels," *arXiv preprint arXiv:1910.00324*, 2019.
- [41] C. Wang, S. Pan, R. Hu, G. Long, J. Jiang, and C. Zhang, "Attributed graph clustering: A deep attentional embedding approach," in *International Joint Conference on Artificial Intelligence*, pp. 3670–3676, 2019.
- [42] G. A. Miller, "Wordnet: a lexical database for english," *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, 1995.
- [43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.



**Xiyao Liu** received the B.S. degree in telecommunication engineering from Tianjin University, Tianjin, China, in 2015. She is currently pursuing a Ph.D. degree in the School of Electrical and Information Engineering, Tianjin University. Her research interests include few-shot learning, human-object interaction, and computer vision.



**Zhong Ji** received the Ph.D. degree in signal and information processing from Tianjin University, Tianjin, China, in 2008. He is currently an Associate Professor with the School of Electrical and Information Engineering, Tianjin University. He has authored over 80 scientific papers. His current research interests include multimedia understanding, zero/few shot learning, cross-modal analysis, and video summarization.



**Yanwei Pang** received the Ph.D. degree in electronic engineering from the University of Science and Technology of China, Hefei, China, in 2004. He is currently a Professor with the School of Electrical and Information Engineering, Tianjin University, Tianjin, China. He has authored over 120 scientific papers. His current research interests include object detection and recognition, vision in bad weather, and computer vision.



**Jungong Han** is currently a tenured Associate Professor of data science with the University of Warwick, U.K. His research interests include computer vision, artificial intelligence, and machine learning. He has published more than 70 articles in top venues, such as IEEE/ACM Transactions (40) and A\* conference (40+, ICML 1; CVPR 3; ICCV 4; ACMMM 3, etc). He is an Associate Editor of Elsevier Neurocomputing, an Editorial Board Member of Springer Multimedia Tools and Applications, and an Associate Editor of IET Computer Vision.



**Xuelong Li** received the Ph.D. degree from the University of Science and Technology of China, Hefei, China, in 2002. He is a Full Professor with the School of Computer Science and Center for OPTical IMagery Analysis and Learning (OPTIMAL), Northwestern Polytechnical University, Xi'an, China. He was with the Xi'an Institute of Optics and Precision Mechanics, Chinese Academy of Sciences, Xi'an.